

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

EV355228691

APPLICATION FOR LETTERS PATENT

**Methods and Systems for Maintaining an Encrypted
Video Memory Subsystem**

Inventor(s):
Glenn F. Evans
Paul England
Nick Wilt

ATTORNEY'S DOCKET NO. ms1-1345usc1

RELATED APPLICATIONS

This application is a continuation application of, and claims priority to, U.S. Patent Application Serial No. 10/314,896, filed on December 9, 2002, the disclosure of which is incorporated by reference herein. This application is also related to the following U.S. Patent Applications, the disclosures of which are incorporated by reference herein:

- U.S. Patent Application Serial No. 10/052,840, entitled "Secure Video Card Methods and Systems," filed January 16, 2002;
- U.S. Patent Application Serial No. 10/178,822, entitled "Methods and Systems Providing Per Pixel Security and Functionality," filed June 24, 2002; and
- U.S. Patent Application Serial No. 10/178,804, entitled "Systems and Methods for Securing Video Card Output," filed June 24, 2002.

TECHNICAL FIELD

This invention relates to methods and systems for processing data using video cards.

BACKGROUND

Typically, a content author, such as a movie studio or a user publishing content on the web, will publish video content that has restrictions on how users can view it. This content can typically be viewed or rendered on a computer such as a personal computer. A great deal of time, effort and money is spent each year by unscrupulous individuals and organizations trying to steal or otherwise inappropriately obtain such video content.

One of the points of attack can be the computer on which such video content is to be viewed or rendered. That is, rogue programs or devices can and often do try to inappropriately obtain video content once it has been received on a computer, such as a personal computer. Among other computer components, this attack can be waged against the video card that processes the video content and/or the bus that transports the video content to and from the video card.

Fig. 1 shows an exemplary video (or graphics) card 100 that includes a bus connector 102 that inserts into a port on a typical computer. Video card 100 also includes a monitor connector 104 (e.g. a 15-pin plug) that receives a cable that connects to a monitor. Video card 100 can include a digital video-out socket 106 that can be used for sending video images to LCD and flat panel monitors and the like.

The modern video card consists of four main components: the graphics processor unit (GPU) 108, the video memory 110, the random access memory digital-to-analog converter (RAMDAC) 112, and the driver software which can be included in the Video BIOS 114.

GPU 108 is a dedicated graphics processing chip that controls all aspects of resolution, color depth, and all elements associated with rendering images on the monitor screen. The computer's central processing unit or CPU (not shown) sends a set of drawing instructions and data, which are interpreted by the graphics card's proprietary driver and executed by the card's GPU 108. GPU 108 performs such operations as bitmap transfers and painting, window resizing and repositioning, line drawing, font scaling and polygon drawing. The GPU 108 is designed to handle these tasks in hardware at far greater speeds than the software running on the system's CPU. The GPU then writes the frame data to the frame buffer (or on-

1 board video memory 110). The GPU greatly reduces the workload of the system's
2 CPU.

3 The memory that holds the video image is also referred to as the frame
4 buffer and is usually implemented on the video card itself. In this example, the
5 frame buffer is implemented on the video card in the form of memory 110. Early
6 systems implemented video memory in standard DRAM. However, this requires
7 continual refreshing of the data to prevent it from being lost and cannot be
8 modified during this refresh process. The consequence, particularly at the very
9 fast clock speeds demanded by modern graphics cards, is that performance is
10 badly degraded.

11 An advantage of implementing video memory on the video card itself is
12 that it can be customized for its specific task and, indeed, this has resulted in a
13 proliferation of new memory technologies:

- 14 • Video RAM (VRAM): a special type of dual-ported DRAM, which
15 can be written to and read from at the same time. It also requires far
16 less frequent refreshing than ordinary DRAM and consequently
17 performs much better;
 - 18 • Windows RAM (WRAM): as used by the Matrox Millennium card,
19 is also dual-ported and can run slightly faster than conventional
20 VRAM;
 - 21 • EDO DRAM: which provides a higher bandwidth than DRAM, can
22 be clocked higher than normal DRAM and manages the read/write
23 cycles more efficiently;
- 24
25

SDRAM: Similar to EDO RAM except the memory and graphics chips run on a common clock used to latch data, allowing SDRAM to run faster than regular EDO RAM;

- SGRAM: Same as SDRAM but also supports block writes and write-per-bit, which yield better performance on graphics chips that support these enhanced features; and
- DRDRAM: Direct RDRAM is a totally new, general-purpose memory architecture which promises a 20-fold performance improvement over conventional DRAM.

Some designs integrate the graphics circuitry into the motherboard itself and use a portion of the system's RAM for the frame buffer. This is called "unified memory architecture" and is used for reasons of cost reduction only and can lead to inferior graphics performance.

The information in the video memory frame buffer is an image of what appears on the screen, stored as a digital bitmap. But while the video memory contains digital information its output medium -- the monitor -- may use analog signals. The analog signals require more than just an "on" or "off" signal, as it is used to determine where, when and with what intensity the electron guns should be fired as they scan across and down the front of the monitor. This is where RAMDAC 112 comes into play as described below. Some RAMDACs also support digital video interface (DVI) outputs for digital displays such as LCD monitors. In such configurations, the RAMDAC converts the internal digital representation into a form understandable by the digital display.

1 The RAMDAC plays the roll of a “display converter” since it converts the
2 internal digital data into a form that is understood by the display.

3 Even though the total amount of video memory installed on the video card
4 may not be needed for a particular resolution, the extra memory is often used for
5 caching information for the GPU 108. For example, the caching of commonly
6 used graphical items - such as text fonts and icons or images- avoids the need for
7 the graphics subsystem to load these each time a new letter is written or an icon is
8 moved and thereby improves performance. Cached images can be used to queue
9 up sequences of images to be presented by the GPU, thereby freeing up the CPU
10 to perform other tasks.

11 Many times per second, RAMDAC 112 reads the contents of the video
12 memory, converts it into a signal, and sends it over the video cable to the monitor.
13 For analog displays, there is typically one Digital-to-Analog Converter (DAC) for
14 each of the three primary colors the CRT uses to create a complete spectrum of
15 colors. For digital displays, the RAMDAC outputs a single RGB data stream to be
16 interpreted and displayed by the output device. The intended result is the right
17 mix needed to create the color of a single pixel. The rate at which RAMDAC 112
18 can convert the information, and the design of GPU 108 itself, dictates the range
19 of refresh rates that the graphics card can support. The RAMDAC 112 also
20 dictates the number of colors available in a given resolution, depending on its
21 internal architecture.

22 The bus connector 102 can support one or more busses that are used to
23 connect with the video card. For example, an Accelerated Graphics Port (AGP)
24 bus can enable the video card to directly access system memory. Direct memory
25 access helps to make the peak bandwidth many times higher than the Peripheral

1 Component Interconnect (PCI) bus. This can allow the system's CPU to do other
2 tasks while the GPU on the video card accesses system memory.

3 During operation, the data contained in the on-board video memory can be
4 provided into the computer's system memory and can be managed as if it were
5 part of the system's memory. This includes such things as virtual memory
6 management techniques that the computer's memory manager employs. Further,
7 when the data contained in the system's memory is needed for a graphics
8 operation on the video card, the data can be sent over a bus (such as a PCI or AGP
9 bus) to the video card and stored in the on-board video memory 110. There, the
10 data can be accessed and manipulated by GPU 108 as described above.

11 This invention arose out of concerns associated with providing methods and
12 systems for protecting data that is used in connection with a video card.

13 14 **SUMMARY**

15 Methods and systems protect digital content such as premium content like
16 movies, programs, and other types of digital audio/visual content. In some
17 embodiments, an architecture and related methods protect content by maintaining
18 the content in encrypted form, whether the content resides in video card memory
19 (referred to herein as "VRAM"), or some other local or remote memory
20 subsystem. The methods and systems enable video card co-processors, such as the
21 graphics processing unit (GPU) to manipulate the encrypted content or data. In
22 various embodiments, the content is maintained in an encrypted format and is
23 unencrypted only when the GPU operates upon the data. After the GPU operates
24 upon the data, the resultant data is re-encrypted and written to memory.

BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 is a block diagram that shows various components of an exemplary video or graphics card that is intended for use in a computer system.

Fig. 2 is a block diagram of an exemplary computer system that can employ video cards in accordance with the described embodiment.

Fig. 3 is a block diagram that shows various components of an exemplary video or graphics card in accordance with one embodiment.

Fig. 4 is a flow diagram that describes steps in a method in accordance with one embodiment.

Fig. 5 is a block diagram that shows various components of an exemplary video or graphics card that is intended for use in a computer system, in accordance with one embodiment.

Fig. 6 is a block diagram that shows various components that can be utilized to implement one or more embodiments.

DETAILED DESCRIPTION

Exemplary Computer System

Fig. 2 illustrates an example of a suitable computing environment 200 on which the system and related methods described below can be implemented.

It is to be appreciated that computing environment 200 is only one example of a suitable computing environment and is not intended to suggest any limitation as to the scope of use or functionality of the media processing system. Neither should the computing environment 200 be interpreted as having any dependency or requirement relating to any one or combination of components illustrated in the exemplary computing environment 200.

1 The various described embodiments can be operational with numerous
2 other general purpose or special purpose computing system environments or
3 configurations. Examples of well known computing systems, environments,
4 and/or configurations that may be suitable for use with the media processing
5 system include, but are not limited to, personal computers, server computers, thin
6 clients, thick clients, hand-held or laptop devices, multiprocessor systems,
7 microprocessor-based systems, set top boxes, programmable consumer electronics,
8 network PCs, minicomputers, mainframe computers, distributed computing
9 environments that include any of the above systems or devices, and the like.

10 In certain implementations, the system and related methods may well be
11 described in the general context of computer-executable instructions, such as
12 program modules, being executed by a computer. Generally, program modules
13 include routines, programs, objects, components, data structures, etc. that perform
14 particular tasks or implement particular abstract data types. The embodiments can
15 also be practiced in distributed computing environments where tasks are
16 performed by remote processing devices that are linked through a communications
17 network. In a distributed computing environment, program modules may be
18 located in both local and remote computer storage media including memory
19 storage devices.

20 In accordance with the illustrated example embodiment of Fig. 2,
21 computing system 200 is shown comprising one or more processors or processing
22 units 202, a system memory 204, and a bus 206 that couples various system
23 components including the system memory 204 to the processor 202.

24 Bus 206 is intended to represent one or more of any of several types of bus
25 structures, including a memory bus or memory controller, a peripheral bus, an

1 accelerated graphics port, and a processor or local bus using any of a variety of
2 bus architectures. By way of example, and not limitation, such architectures
3 include Industry Standard Architecture (ISA) bus, Micro Channel Architecture
4 (MCA) bus, Enhanced ISA (EISA) bus, Video Electronics Standards Association
5 (VESA) local bus, and Peripheral Component Interconnects (PCI) bus also known
6 as Mezzanine bus.

7 Computer 200 typically includes a variety of computer readable media.
8 Such media may be any available media that is locally and/or remotely accessible
9 by computer 200, and it includes both volatile and non-volatile media, removable
10 and non-removable media.

11 In Fig. 2, the system memory 204 includes computer readable media in the
12 form of volatile, such as random access memory (RAM) 210, and/or non-volatile
13 memory, such as read only memory (ROM) 208. A basic input/output system
14 (BIOS) 212, containing the basic routines that help to transfer information
15 between elements within computer 200, such as during start-up, is stored in ROM
16 208. RAM 210 typically contains data and/or program modules that are
17 immediately accessible to and/or presently be operated on by processing unit(s)
18 202.

19 Computer 200 may further include other removable/non-removable,
20 volatile/non-volatile computer storage media. By way of example only, Fig. 2
21 illustrates a hard disk drive 228 for reading from and writing to a non-removable,
22 non-volatile magnetic media (not shown and typically called a "hard drive"), a
23 magnetic disk drive 230 for reading from and writing to a removable, non-volatile
24 magnetic disk 232 (e.g., a "floppy disk"), and an optical disk drive 234 for reading
25 from or writing to a removable, non-volatile optical disk 236 such as a CD-ROM,

1 DVD-ROM or other optical media. The hard disk drive 228, magnetic disk drive
2 230, and optical disk drive 234 are each connected to bus 206 by one or more
3 interfaces 226.

4 The drives and their associated computer-readable media provide
5 nonvolatile storage of computer readable instructions, data structures, program
6 modules, and other data for computer 200. Although the exemplary environment
7 described herein employs a hard disk 228, a removable magnetic disk 232 and a
8 removable optical disk 236, it should be appreciated by those skilled in the art that
9 other types of computer readable media which can store data that is accessible by a
10 computer, such as magnetic cassettes, flash memory cards, digital video disks,
11 random access memories (RAMs), read only memories (ROM), and the like, may
12 also be used in the exemplary operating environment.

13 A number of program modules may be stored on the hard disk 228,
14 magnetic disk 232, optical disk 236, ROM 208, or RAM 210, including, by way of
15 example, and not limitation, an operating system 214, one or more application
16 programs 216 (e.g., multimedia application program 224), other program modules
17 218, and program data 220. A user may enter commands and information into
18 computer 200 through input devices such as keyboard 238 and pointing device 240
19 (such as a "mouse"). Other input devices may include a audio/video input
20 device(s) 253, a microphone, joystick, game pad, satellite dish, serial port, scanner,
21 or the like (not shown). These and other input devices are connected to the
22 processing unit(s) 202 through input interface(s) 242 that is coupled to bus 206,
23 but may be connected by other interface and bus structures, such as a parallel port,
24 game port, or a universal serial bus (USB).

1 A monitor 256 or other type of display device is also connected to bus 206
2 via an interface, such as a video adapter or video/graphics card 244. In addition to
3 the monitor, personal computers typically include other peripheral output devices
4 (not shown), such as speakers and printers, which may be connected through
5 output peripheral interface 246.

6 Computer 200 may operate in a networked environment using logical
7 connections to one or more remote computers, such as a remote computer 250.
8 Remote computer 250 may include many or all of the elements and features
9 described herein relative to computer.

10 As shown in Fig. 2, computing system 200 is communicatively coupled to
11 remote devices (e.g., remote computer 250) through a local area network (LAN)
12 251 and a general wide area network (WAN) 252. Such networking environments
13 are commonplace in offices, enterprise-wide computer networks, intranets, and the
14 Internet.

15 When used in a LAN networking environment, the computer 200 is
16 connected to LAN 251 through a suitable network interface or adapter 248. When
17 used in a WAN networking environment, the computer 200 typically includes a
18 modem 254 or other means for establishing communications over the WAN 252.
19 The modem 254, which may be internal or external, may be connected to the
20 system bus 206 via the user input interface 242, or other appropriate mechanism.

21 In a networked environment, program modules depicted relative to the
22 personal computer 200, or portions thereof, may be stored in a remote memory
23 storage device. By way of example, and not limitation, Fig. 2 illustrates remote
24 application programs 216 as residing on a memory device of remote computer
25 250. It will be appreciated that the network connections shown and described are

1 exemplary and other means of establishing a communications link between the
2 computers may be used.

3 4 Overview

5 The various methods and systems described herein are directed to
6 protecting content such as premium content like movies, programs, and other
7 types of digital audio/visual content. In the described embodiments, an
8 architecture and related methods protect content by maintaining the content in
9 encrypted form, whether the content resides in video card memory (referred to
10 herein as “VRAM”), or some other local or remote memory subsystem. The
11 methods and systems enable video card co-processors, such as the graphics
12 processing unit (GPU) to manipulate the encrypted data. Processing on the GPU
13 can be controlled, in some embodiments, by an application that need not
14 necessarily be entirely trusted.

15 Maintaining the content in encrypted form in the memory is more resistant
16 to security leaks since only select portions of the GPU are able to access
17 unencrypted content.

18 As an overview to an exemplary system, consider the Fig. 3 system at 300
19 which represents some of the components that can reside on a video card. There,
20 system 300 comprises a graphics processor unit 302 having multiple inputs 302a,
21 302b, and an output 302c. In this example and for simplicity, only two inputs are
22 illustrated for the GPU. Typically, however, GPU’s have more than two inputs.
23 In some embodiments, the GPU has eight inputs. Each of the GPU’s inputs is
24 associated with a portion of memory (also referred to as a “surface”) that holds
25 data that is subject to processing by the GPU. In this example, since there are two

1 GPU inputs, there are two surfaces 304, 306. Thus, the pixels from surface 304
2 constitute one input to the GPU and the pixels from surface 306 constitute the
3 other input for the GPU.

4 The GPU is configured to operate upon the data of each of surfaces 304,
5 306 and provide the output of its operation onto an output surface 308. This
6 output surface can then be read and rendered to a display by other components that
7 reside on the video card. The operations that can be performed by the GPU are
8 also termed "programs" and can be represented as mathematical operations such
9 as additions, subtractions, multiplications and the like, or include control
10 instructions such as looping or branching, as will be appreciated and understood
11 by the skilled artisan. That is, the GPU might take values associated with
12 individual pixels of surface 304, and add those values to values associated with
13 individual pixels of surface 306, and write the corresponding result to a pixel
14 address on surface 308.

15 Notice in the illustration that the data on each of surfaces 304, 306 and 308
16 is encrypted. In accordance with one embodiment, a cryptographic processor 310
17 is associated with the video card and represents a trusted component. In this
18 example, the cryptographic processor comprises a hardware component in the
19 form of an integrated circuit chip that is physically mounted on the video card.
20 The cryptographic processor is responsible for setting up decryptors and
21 encryptors to assist in decryption and encryption operations on the video card.
22 Accordingly, decryptors 312, 314 and encryptor 316 are associated with
23 cryptographic processor 310. In this illustration, the decryptors and encryptor are
24 shown to logically reside between the GPU 302 and a surface of the VRAM. The
25 decryptors and encryptors can be physically located in other places. For example,

1 the GPU can have specially configured encryption hardware which is configurable
2 by the cryptographic processor, as noted below.

3 In one embodiment, each individual surface of the VRAM that is to hold
4 encrypted content is associated with its own encryptor/decryptor. The
5 encryption/decryption algorithms and keys that are associated with a particular
6 surface can be unique for that surface. Thus, decryptor 312 uses a unique key
7 associated with data on surface 304 to decrypt the data and enable the GPU to
8 process the data. Similarly, decryptor 314 uses a unique key associated with data
9 on surface 306 to decrypt the data and enable the GPU to process the data. The
10 output surface 308 to which the resultant data is written by the GPU has its own
11 associated encryptor which encrypts the resultant data and writes the encrypted
12 data to the surface. The index of the key within the cryptoprocessor can be
13 returned to the application to be used to identify (by associating the appropriate
14 surface with it) which keys should be used within the decryption and re-encryption
15 operations.

16 In this example, data that resides in the VRAM (or some other local or
17 remote memory) is always kept in encrypted form. The data is only decrypted
18 when the GPU is to operate upon it and then after the operation, before the data is
19 written to the VRAM, it is re-encrypted.

20 In this example, whenever the GPU 302 wishes to perform an operation on
21 encrypted data that resides on a surface, on a per pixel basis, a surface-associated
22 decoder, under the influence of cryptoprocessor 310, decrypts each pixel and
23 provides the pixels to the GPU for the operation. After the operation, an
24 encryptor, under the influence of the cryptoprocessor 310, re-encrypts the output
25 of the GPU to provide the encrypted result to another surface of the VRAM.

1 As an added measure of safety, the various operations or programs that can
2 be performed by the GPU can be restricted if encrypted output is available. That
3 is, it is possible for some GPU operations to permit resultant data to be written to a
4 cache and then later written out to an external memory location without being
5 encrypted. In these situations, the video card can be programmed to disallow
6 those types of operations thus ensuring the protection of the data.

7 This way, the only time that data is actually in the open in an unencrypted
8 form is when the pixel inputs are being provided to or within, and operated upon
9 by the GPU. Now, once the resultant data has been provided onto a different
10 surface, for example surface 308, the DAC can read the encrypted surface, decrypt
11 it and then display the data to a suitable display.

12 An added benefit of this system is that the GPU is actually allowed to
13 perform operations on the data which can greatly accelerate the graphics
14 processing capabilities of the video card. This system also maintains the video
15 memory in protected, encrypted form so that if any components, rogue programs
16 and the like start snooping around the VRAM, all that is present is encrypted data.

17 Fig. 4 is a flow diagram that describes steps in a method in accordance with
18 one embodiment. The method can be implemented in connection with any
19 suitable hardware, software, firmware or combination thereof. In the illustrated
20 and described embodiment, the method can be implemented in connection with the
21 systems described above and below.

22 Step 400 receives encrypted data. The encrypted data can typically
23 comprise some form of protected content. Step 402 writes the encrypted data to
24 memory. This memory can comprise local or remote memory. In the examples
25 given above, the memory comprises the video card's memory such as the VRAM.

1 Step 404 decrypts the encrypted data. This step can be performed responsive to an
2 indication that the data is to undergo some type of operation by the GPU. This
3 indication can come from an application which effectively notifies the video card
4 to perform some type of operation on data that it specifies. It is not necessary for
5 the application to be protected or trusted as it is not necessary for the application
6 to access the data. In addition, this step can be performed by a suitably configured
7 decryption component. In the Fig. 3 example, such a component is provided in the
8 form of a decryptor that is associated with a hardware cryptographic processor that
9 controls encryption and decryption capabilities on the video card. Further, in the
10 Fig. 3 example, each surface or memory portion that is to hold encrypted data has
11 its own associated encryptor and decryptor. It is to be noted and appreciated that
12 while the decryptors and encryptor that are illustrated in Fig. 3 are shown to reside
13 outside the GPU, such need not be the case. That is, the GPU can have its own
14 cryptographic hardware that is communicatively associated with the cryptographic
15 processor. In this situation, the decryption and encryption takes place inside the
16 GPU, and not externally as might be suggested by the Fig. 3 illustration. The keys
17 are controlled by the cryptographic processor and are not available to other
18 components within the GPU (other than the encryptors and decryptors).

19 Step 406 operates on the decrypted data with the GPU to provide resultant
20 data. Any suitable GPU operations can be performed, as noted above. Step 408
21 re-encrypts the resultant data. This step can be implemented by a suitably
22 configured encryption component, an example of which is provided above. Step
23 410 writes the encrypted data to memory. This step can be implemented by
24 writing the encrypted data to any suitable memory. In the illustrated and described
25 example, the encrypted data is written to a VRAM surface that is compatible with

1 the surfaces from which the data was originally read. Step 412 decrypts and
2 displays the encrypted data. This step can be performed by a suitably configured
3 display converter such as a RAMDAC.

4 Thus, the above system can maintain data in the VRAM in encrypted form,
5 on a per-pixel basis, anytime when a GPU operation is not being performed on the
6 data. Additionally, the encrypted data can be decrypted on a pixel-by-pixel basis
7 and accordingly processed by the GPU before being re-encrypted and written back
8 out to the VRAM.

9 It is to be appreciated and understood that any number of suitable
10 encryption/decryption paradigms could be utilized in connection with and to
11 implement the above-described system without departing from the spirit and scope
12 of the claimed subject matter.

13 14 Exemplary Architecture

15 Fig. 5 shows an exemplary video (or graphics) card 500 architecture in
16 accordance with one embodiment. Card 500 includes a bus connector 502 that
17 plugs into a port on a typical computer. Video card 500 also includes a monitor
18 connector 504 (e.g. a 15-pin plug) that receives a cable that connects to a monitor.
19 Video card 500 can, but need not, include a digital video-out (e.g. DVI) socket 506
20 that can be used for sending video images to digital displays and the like.

21 Like the video card of Fig. 1, video card 500 comprises a graphics
22 processor unit (GPU) 508, video memory 510, display converter or random access
23 memory digital-to-analog converter (RAMDAC) 512, and driver software which
24 can be included in the Video BIOS 514.

GPU 508 is a dedicated graphics processing chip that controls all aspects of resolution, color depth, and all elements associated with rendering images on the monitor screen. The memory controller (sometimes integrated into the GPU) manages the memory on the video card. The computer's central processing unit or CPU (not shown) sends a set of drawing instructions and data, which are interpreted by the graphics card's proprietary driver and executed by the card's GPU 508. GPU 508 performs such operations as bitmap transfers and painting, window resizing and repositioning, line drawing, font scaling and polygon drawing. The GPU can then write the frame data to the frame buffer (or on-board video memory 310). In the illustrated and described embodiment, GPU 508 can comprise cryptographic hardware 508a which can assist in cryptography, as described in more detail below.

The information in the video memory frame buffer is an image of what appears on the screen, stored as a digital bitmap. RAMDAC 512 is utilized to convert the digital bitmap into a form that can be used for rendering on the monitor, as described above.

In addition to these components, in this embodiment, video card 500 comprises a memory controller 516 which can include a cache (not specifically illustrated), a cryptographic processor 518 that can include a key manager 520, as well as a bank of keys 522. Although illustrated as part of the cryptographic processor, the key manager 520 can comprise a separate component.

Memory controller 516 receives data on the video card and manages the data in the video memory 510. The memory controller can also be responsible for managing data transfers between the video card and system memory.

1 Cryptographic processor 518 is responsible for organizing cryptographic
2 functions that take place on the video card.

3 It is desirable for secure graphics cards, such as card 500, to be able to
4 authenticate themselves as such. In particular, it is desirable for trusted software,
5 such as secure application 524, to be able to distinguish a secure graphics card
6 from a traditional graphics card or a circumvention device. In addition, it is
7 desirable for trusted software to be able to reveal cryptographic keys to the
8 graphics card and to be able to verify that the receiver of the keys is indeed a
9 secure graphics card. For this purpose, in accordance with the described
10 embodiment, secure graphics cards such as card 500 are equipped with
11 cryptographic processor 518, which performs standard cryptographic tasks of
12 authentication and key transport.

13 In accordance with the described embodiment, cryptographic processor 518
14 is individualized and certified during manufacture. Individual cryptographic
15 processors can contain a unique private decryption key K_{priv} . Although subject to
16 change depending on different requirements and design constraints, the associated
17 encryption/decryption algorithm can be RSA, and the key length can be 1024 bits.
18 The cryptographic processor can be permanently attached to graphics card 500,
19 either by adding it to an existing chip or by adding it as a separate chip to the card.

20 In the illustrated and described embodiment, cryptographic processor 518
21 can implement a public key crypto algorithm (as defined below) and hides a
22 unique private key. It can perform one public key decryption and can utilize a
23 public key accelerator. In addition, the cryptographic processor can implement a
24 symmetric cipher (AES) and some control logic.
25

1 In one embodiment, the cryptographic processor has the following volatile
2 registers.

- 3
- 4 • A 256-bit register S for the session key. The lifetime of this key is
5 typically the running time of the trusted software.
- 6 • An array of x (x TBD) index keys. Each key is 128 bits long. Each key
7 can be associated with a particular surface and can be used by the
8 graphics card to decrypt its contents. The lifetime of each key can be
9 governed by instructions from the trusted software.

10 As noted above, in the illustrated and described embodiment, the
11 cryptographic processor is permanently attached to the graphics card. In this
12 example, there are two interfaces to cryptographic processor 518--an external
13 interface to trusted software 524, and an interface to GPU 508. In the illustrated
14 and described embodiment, the interface to the trusted software 524 is
15 standardized, while the interface to the GPU 508 can be implementation-specific.

16 17 The External Interface

18 The external interface can use the basic PK encryption protocol for
19 authentication and key transport. Under this protocol, trusted software 524
20 encrypts a session key with the public key of cryptographic processor 518. The
21 cryptographic processor receives the resulting cryptoblob and decrypts it with its
22 private key, thus obtaining the session key. Now, the trusted software and the
23 cryptographic processor share a secret. The trusted software can use this session
24 key to send instructions to the cryptographic processor. At an abstract level, the
25

1 external interface can be exposed through various functions by the cryptographic
2 processor.

3 4 The Internal Interface

5
6 The term “internal interface” refers to the interface between cryptographic
7 processor 518 and the rest of graphics card 500. The cryptographic processor can
8 use this interface to configure the GPU’s cryptographic hardware. In one
9 embodiment, the details of this interface are up to the implementation of each
10 individual graphics card, subject to the following restrictions:

- 11
12 • Removal of the cryptographic processor from the graphics card
13 should not be trivial. If the cryptographic processor is implemented
14 as a separate chip, this is mainly a restriction on the mechanical
15 interface, which attaches the cryptographic processor to the graphics
16 card. Typically, the cryptographic processor should be soldered onto
17 the graphics card. Alternatively, the cryptographic processor could
18 reside on the same chip as the main GPU. Use of standardized
19 mechanical interfaces that allow the cryptographic processor to be
20 removed (for example, a socket-mounted smart card reader) is not
21 desirable.
- 22 • The physical connection between the cryptographic processor and
23 the rest of the graphics card should not be accessible and should not
24 be exposed through standard interfaces. For example, a USB
25 connector on this bus is not desirable.

21 22 Implementation Example

23 Fig. 6 illustrates diagrammatically but one particular implementation
24 example. This illustration is not intended to limit application of the claimed
25 subject matter. Rather, such illustration is intended to illustrate but one way in

1 which the various inventive features described in this document can be
2 implemented. In this example, cryptographic management is implemented by a
3 discrete cryptographic processor which interfaces with the GPU on pixel accesses.

4 In this particular example, components of the video card include
5 cryptoprocessor 600 which serves as a key repository and key distributor, GPU
6 602, video memory or VRAM 604 and a DAC/DVI component 606 that is
7 configured to display data to a suitable display device. Inside the GPU 602,
8 various components can include so-called pixel shaders that effectively perform
9 programs or operations on the data that it receives. Notice in this example that the
10 GPU comprises encryption and decryption capabilities (e.g. encryption/decryption
11 hardware) as indicated by encryptor/decryptor 602a and decryptor 602b. Here, the
12 decryptor 602b is added to the GPU's texture mapping unit on the input side, and
13 encryptor/decryptor 602a is added to the alpha blending unit on the output side. In
14 implementing this particular functionality, hardware designers can follow some
15 rules to facilitate implementation. Specifically, since stream ciphers do not enable
16 random access to encrypted data, the system should use block ciphers, e.g.
17 encrypting the data 128-bits at a time. The texture mapping unit can decrypt on a
18 cache line fill, and the alpha blending unit can decrypt when reading a cache line
19 from the color buffer and encrypt before writing. The encryption keys used in
20 these operations can, and should often be different. Computational tasks other
21 than 3D rendering, such as video decoding, are straightforward extensions of the
22 just-described paradigm. Instead of textures, video macroblocks would serve as
23 the encrypted input; instead of a color buffer, the output frame being decoded
24 would serve as the encrypted output. If content must be protected as it is delivered
25

1 inband in the command stream to the GPU, the command buffer may be encrypted
2 as well.

3 Video memory 604 comprises multiple surfaces 604a, 604b which serve to
4 hold encrypted content which, in this example, is represented by encrypted
5 premium content 610. Video memory 604 also includes one or more resultant
6 surfaces 604c, and a desktop or primary surface 604d that is read, decrypted and
7 output by DAC/DVI component 606.

8 This system also includes an application 608 that has a trusted portion that
9 sets up the encryption capabilities of the cryptographic processor 600. That is,
10 application 608 has a trusted component that will set up the keys with the
11 cryptographic processor 600. The application really does not need access to the
12 data. Rather, the application can really simply direct GPU 602 on what operations
13 should be performed. This enables the application to leverage the GPU so that the
14 GPU itself is responsible for manipulating the data.

15 For example, assume that application 608 wants to do picture-in-picture
16 (PIP) video. The application can create three encrypted surfaces, but does not
17 need to know anything about the video. When the application creates the surfaces
18 (in VRAM), it communicates with the cryptographic processor 600 and creates or
19 otherwise indicates or provides a key associated with each of the surfaces to the
20 cryptographic processor. The key can be managed in a bank of keys in the
21 cryptographic processor 600. In this embodiment and as noted above, there is one
22 key per associated encrypted surface.

23 Now, when operating system tells the graphics card to run a particular
24 program on these created surfaces, the graphics driver will identify the input
25 surfaces (such as surfaces 604a, 604b), and will then ask the cryptographic

1 processor 600 to select keys associated with the surfaces of interest which, in turn,
2 initializes the encryptors/decryptors in the GPU.

3 In one implementation, the encryption and decryption keys (or sets of keys)
4 can be correlated by their index in the cryptoprocessor. For example, the trusted
5 portion of the application of 608 will negotiate the keys with the cryptoprocessor
6 600 and it will return the key indices to the untrusted portion of the application
7 608. The untrusted application will create surfaces which are identified as being
8 encrypted or decrypted using the key indices. When the GPU 602 performs
9 operations, it will use the key indices to initialize the encryption and decryption
10 keys on the surfaces. The cryptoprocessor will receive the request from the GPU
11 to set the indices on the encryptors and decryptors. The cryptoprocessor can
12 decide to successfully setup the encryptors (304 and 306 in Fig 3) and decryptors
13 (316 in fig 3) if the indices are compatible (i.e. it is valid to transcript from each
14 input key to the output key). The cryptoprocessor can map each key index into the
15 actual key to transfer to the encryptor and decryptor component. If the untrusted
16 application attempts to setup an invalid configuration, then the operation will fail
17 or could produce unusable data (i.e. go ahead with the decryption with invalid
18 keys producing incorrectly decrypted data).

19 When GPU 602 actually runs the program on the pixels, decryptor 602b
20 reads a pixel, decrypts it with the appropriate key, the GPU operates on the pixel
21 with the program, and then the GPU re-encrypts the resultant pixel (as with
22 encryptor/decryptor 602a) and writes it to an appropriate surface, such as surface
23 604c in VRAM 604. The DAC/DVI component 606 can then decrypt the data for
24 subsequent display.

1 In the past, the processing of the video data would typically be performed
2 off of the video card by the computer system's central processing unit (CPU)
3 under the influence of the application. This approach to video processing does not
4 scale and, while adequate for fairly small pictures, is entirely inadequate for larger
5 pictures such as those employed in HDTV. In this past scenario, the GPU's role
6 was really relegated to that of the DAC functionality on the back end. In order to
7 make video processing algorithms scale, hardware acceleration is needed. By
8 performing the operations using the GPU, one can achieve the hardware
9 acceleration that is necessary to provide desirable, scalable performance. In
10 addition, while providing the desired scalability, the above-described systems and
11 methods can be utilized to protect the data any time when it is not being processed
12 by the GPU. As the data is maintained in encrypted form, the data is fairly well
13 protected against theft and use by unauthorized components or parties.

14 As a further protective measure and to protect the content from being
15 moved on the desktop surface 604d, the content can be encrypted on the desktop
16 surface by either including position information or transcribing to a position
17 dependent global key. The position can also be enforced by limiting which
18 processes can update the origin (or clipping lists) and the key table. With a steam
19 cipher, the second encryption can be applied before the first encryption is removed
20 ensuring that the clear stream is not visible.

21 22 **Memory Optimizations**

23 The following section describes various optimizations that can be utilized
24 in connection with the encryption and decryption that takes place on the video
25 card.

1 Since textures and off-screen surfaces typically require random access, it is
2 advantageous that they be encoded with block ciphers. There is good synergy
3 between the typical block size for a block cipher and the typical cache line size for
4 a modern 3D accelerator. That is, if the cache line and block size are both 128 bits
5 (or bear some integer-size relation), then efficient encryption and decryption can
6 be implemented in the hardware. Even if there are slight differences (for example,
7 block size of 128 bits and cache line size of 256 bits) the hardware implementation
8 can be efficient.

9 The key size and usage pattern must be chosen to match the cache structure
10 of the video accelerator to feasibly implement the per pixel encryption. In fact, the
11 decryptors and encryptors can be moved logically into the memory controller
12 portion of the GPU. Consider the model where pixel data is read into a cache page
13 on a pixel access (or copied from the cache if it is already present) and pixel data
14 is written to the write cache then copied to memory on a write cache page eviction.
15 Instead of decrypting on a per pixel basis, the memory can be decrypted on a per
16 cache page read. If the pixel data has already been decrypted (the key index could
17 be used to decide this), then the decryption can be avoided. That is, the GPU's
18 memory controller can optimize decryption and re-encryption iterations by
19 caching decrypted pages in a local page pool cache. If a page is not present in the
20 cache, then it can be decrypted and placed into the cache. When a page is flushed
21 and written out to memory, it is first re-encrypted then written. Hence decryption
22 can be avoided if the same page needs to be accessed many times. Since there
23 could be an integer relationship between the tile/block-key/page-size, then this
24
25

1 algorithm should greatly exploit coherency and greatly reduce decryption and re-
2 encryption accesses to memory.

3 Output data only needs to be re-encrypted when a page write occurs. The
4 cache needs to be purged of encrypted data before and after the GPU program has
5 completed.

6 Another optimization is to maintain the cache between programs if the
7 same output key is used. However, during the period between operations, the rest
8 of the GPU should be blocked from accessing cache pages containing decrypted
9 content. This optimization would be necessary for performing operations such as
10 video decoding which perform hundreds of operations on the same input and
11 output data.

12 **Implementing Encryption on current Swizzled Hardware designs**

13 One problem with encrypted texture data is that a block encryption scheme
14 requires an adjacent block of bytes to be available before it can be encrypted or
15 decrypted; and a cache line fill requires that the pixel data be “swizzled”—that is,
16 the translation from (X,Y) position in the image to an address must be formed
17 such that the cache line fill yields a 2D region of pixels. To date, hardware
18 vendors have exposed ostensibly linear surface formats while swizzling image
19 data without the knowledge of the application. Since trusted software will be
20 emitting the encrypted texture data, however, it must have *a priori* knowledge of
21 the swizzling scheme so it can encrypt adjacent blocks of data and preserve 2D
22 locality. A good solution is to define a dictionary of swizzled-image formats
23 (including YUV 4:4:4, 4:2:2, and 4:2:0 as well as RGB formats) for use by the
24 application. The performance of these formats may not be quite as high as if the
25

1 images were swizzled to a hardware-specific format, but the encryption is
2 presumably worth a slight performance degradation.

3 4 **Conclusion**

5 The various methods and systems described above can protect content such
6 as premium content like movies, programs, and other types of digital audio/visual
7 content. In the described embodiments, an architecture and related methods
8 protect content by maintaining the content in encrypted form, whether the content
9 resides in video card memory (referred to herein as “VRAM”), or some other local
10 or remote memory subsystem. The methods and systems enable video card co-
11 processors, such as the graphics processing unit (GPU) to manipulate the
12 encrypted data. Processing on the GPU can be controlled, in some embodiments,
13 by an application that need not necessarily be entirely trusted. Maintaining the
14 content in encrypted form in the memory is more resistant to security leaks since
15 only select portions of the GPU are able to access unencrypted content.

16 Although the invention has been described in language specific to structural
17 features and/or methodological steps, it is to be understood that the invention
18 defined in the appended claims is not necessarily limited to the specific features or
19 steps described. Rather, the specific features and steps are disclosed as preferred
20 forms of implementing the claimed invention.